

SVM-Based Anomaly Detection for Network Intrusion Systems: An Advanced Approach to Intrusion Prevention

Arihant Singhvi
Reg no: 220962214
Roll no.: 36

Pranav Menon
Reg no: 220962098
Roll no: 28

Abhinav Nambiar
Reg no: 220962242
Roll no.: 38

Abstract—This project presents the development of an Intrusion Detection System (IDS) utilizing Support Vector Machine (SVM) algorithms to enhance network security through anomaly detection, addressing the limitations of traditional signature-based methods that struggle with identifying novel threats. By training an SVM model on labeled network traffic data, the system effectively distinguishes between normal and anomalous patterns, enabling the detection of sophisticated attacks. Comprehensive experimentation demonstrates high detection accuracy with a low false positive rate. The project also tackles challenges such as the need for high-quality labeled data, significant computational resources, and the complexity of data preprocessing, emphasizing the importance of meticulous data selection. Additionally, it reviews current anomaly-based IDS techniques, comparing them with the proposed SVM-based approach, and identifies key challenges like computational constraints and model adaptability. The resulting scalable and adaptive IDS framework advances network security by detecting emerging threats and overcoming traditional detection limitations.

Keywords—Anomaly-Based Network Intrusion Detection System (A-NIDS), SVM, Naïve-bayes

I. INTRODUCTION

With rapid growth of cyber threats, the advent of Intrusion Detection Systems (IDS) has brought intrusion detection methods based on traditional signatures into scrapping. While signature-based IDS are very effective at recognizing known attack patterns, they are limited in their ability to detect previously seen or completely novel attacks because they rely exclusively on pre-defined signatures. This limitation, however, necessitates more advanced and unconventional approaches to detect anomalous behaviour in contrast to normal network operations.

Anomaly based intrusion detection constitutes a major advancement in addressing these limitations. Unlike signature-based systems, which depend on matching their attack signatures to known exploits, anomaly-based IDS attempt to detect deviations from the established patterns of normal network behavior. These systems are continuously monitoring and analyzing the network traffic and can identify unusual, or suspicious patterns or activities that may point to malicious threats, which can be zero-day attacks and other advanced intrusion attempts.

This project offers an Anomaly Based Intrusion Detection System (A-NIDS), using Support Vector Machine (SVM) algorithms, to improve the detection of anomalous network behavior. It is because SVM is very robust to high

dimension data and very effective for creating the best decision boundaries that SVM is a very suitable algorithm for this task. Using the SVM's capacity to address data points with precision enables our IDS to differentiate better between normal and anomalous network traffic. With learnt and generalized normal behavior, the system has been trained on a comprehensive dataset of labelled network traffic. This training allows the IDS to detect deviations from these patterns and provide potential threat warning that might not be recognized as an attack signature.

Through this approach, the SVM-based IDS aims to provide a more dynamic and adaptable solution for network security. The ability of the SVM model to detect previously unknown threats and adapt to evolving attack strategies is a significant advantage over traditional methods. This project includes rigorous experimentation to evaluate the performance of the IDS, focusing on achieving high detection accuracy while minimizing false positives. The goal is to deliver an IDS framework that enhances network protection by effectively identifying and responding to both known and emerging threats, thus addressing the limitations of conventional detection systems.

II. SALIENT FEATURES

- **Anomaly Detection Capability:** Utilizes an anomaly-based detection approach to detect deviation from expressively trained network behavior patterns for detection of previously unknown or novel threats.
- **Support Vector Machine (SVM) Algorithm:** Utilizes the strength of the SVM algorithm in the high dimensionality case for robust classification of network traffic into normal and anomalous classes.
- **Comprehensive Data Collection:** Thoroughly acquiring and labeling network traffic data to form a defined dataset suitable for training an IDS model.
- **Advanced Data Preprocessing:** Includes thorough cleaning and normalization of data, alongside feature selection, to enhance the quality and relevance of the input data for model training.
- **Effective Model Training:** Focuses on training the SVM model with preprocessed data to accurately learn and generalize patterns of normal and anomalous network behavior.

- **Comprehensive Performance Evaluation:** Measures the effectiveness of the IDS through performance metrics such as accuracy, precision, recall, and F1 score, ensuring reliable detection of anomalies and minimal false positive rates.

III. METHODOLOGY

In this research, we employ the CSE-CIC-IDS 2018 dataset, which is systematically generated using a profile-based approach. This dataset offers extensive descriptions of various intrusions and abstract distribution models relevant to applications, protocols, and even lower-level network elements. Its use of profiles allows for flexibility, enabling the dataset to accommodate a wide variety of network protocols and topologies. The dataset includes two primary types of profiles: the B-profile and the M-profile. The B-profile encompasses essential protocol characteristics such as packet size distributions, the number of packets per flow, payload size and patterns, and request time distributions. These attributes make it particularly useful for machine learning applications. In contrast, the M-profile focuses on clearly defining specific attack scenarios, offering precise details to simulate different intrusion behaviors. The dataset encompasses both normal and adversarial network behaviors which could be converted to an Intrusion Detection System (IDS) using the Support Vector Machine (SVM) algorithm. The dataset consists of 80 features with over 1 million entries with no null values. Prominent features include destination port, protocol, flow duration, total forwards packet, total backwards packet and label for the anomaly. The label column describes the type of pattern encountered in the computer system based on the data packet configurations which help us distinguish between anomalous and normal systems. The attributes of pattern type are as follows:

1. Normal Pattern:
 - Benign
2. Attack Pattern:
 - DoS attacks-GoldenEye
 - Dos attacks-Slowloris

After analyzing the dataset, steps for data cleaning and preprocessing are initiated. The ‘Timestamp’ column is dropped from dataset since it only provides the date and time of the network packet reading and does not contribute to the prediction model. The ‘Label’ column is checked for its unique values to know the types of normal and attack patterns and their counts. The dataset is checked for categorical values which are only found in ‘Label’ feature. One-hot encoding is carried out which converts categorical data into a series of binary columns, each representing a unique category to improve predictions. The transformation of ‘Label’ feature is carried out which converts normal pattern type to 0 and attack pattern type to 1. Correlation coefficient of the remaining input features with the ‘Label’ column is displayed which gives a general idea for their contribution in predicting the network pattern type.

The train-test split takes place with 33-67% distribution of training and testing sets respectively. The training and testing set for input features is checked for infinite values which may have occurred due to overflow in exponential, undefined mathematical operations and data processing errors. These infinite values are converted to null values which are then imputed to the mean of the respective feature to which they belong. The feature values are then scaled according to the mean and standard deviation of that feature to normalize their range.

Name	Abbreviation	Formula
Cubic	CU	$\varphi(r) = r^3$
Thin plate splines	TPS	$\varphi(r) = r^2 \log(r)$
Generalized Thin plate splines	GTPS	$\varphi(r) = r^{2m} \log(r), m \in \mathbb{N}$
Inverse quadrics(or Cauchy)	IQ	$\varphi(r) = \frac{1}{c^2 + r^2}$
Multiquadrics	MQ	$\varphi(r) = \sqrt{c^2 + r^2}$
Inverse Multiquadrics	IMQ	$\varphi(r) = \frac{1}{\sqrt{c^2 + r^2}}$
Gaussian RBF	GA	$\varphi(r) = e^{-r^2/c^2}$

Definition 2.2 θ -method, ($0 \leq \theta \leq 1$), is general finite-difference approximation to $\frac{\partial^2 u(x,t)}{\partial x^2}$ given by:

$$\frac{\partial^2 u(x,t)}{\partial x^2} \cong \theta \delta_{2,x} U_{i,j+1} + (1-\theta) \delta_{2,x} U_{i,j},$$

Fig. 1. Different types of RBF functions which can handle non-linear data by mapping it into higher-dimensional space in SVM models.

The SVM model is then developed using the Radial Basis Function (RBF) kernel, specifically Gaussian RBF, selected for its ability to handle the non-linear nature of network traffic data. The model is trained on the labeled dataset with a focus on accurately distinguishing between normal and anomalous traffic. To optimize the model, we employ cross-validation techniques to fine-tune hyperparameters such as the penalty parameter (C) and kernel coefficient (gamma).

The formula for RBF kernels between 2 feature vectors x and x' is given by:

$$K(x, x') = e^{(-\gamma \|x - x'\|^2)}$$

where:

- $\|x - x'\|^2$ represents the squared Euclidean distance between the two feature vectors x and x' .
- γ (gamma) is a parameter that determines the influence of each training example.

The RBF kernel fits the dataset into a higher dimension to produce the linear separation between anomalous and normal pattern type in a less complex fashion. It does not explicitly transform the input data and leads to the formation of hyperplane in the higher dimensional space which comfortably separates the classes. The regularization parameter (C) prevents the risk of overfitting in the SVM model. A higher C value allows fewer misclassifications on the training set and increase chances of overfitting and lower C value will allow some misclassification which will decrease accuracy of model but improve its generalization to unforeseen data. Thus, C acts as a penalty parameter for overfitting and typically different values of C are utilized before finding an optimal parameter. Gamma hyperparameter describes the reach of each training example. Lower gamma values lead to a broader influence, making the model more generalized, reduces overfitting and makes the decision boundary smoother while higher values tend to make decision boundary more sensitive to incoming data points and increases chance of overfitting the model. Thus, gamma values are optimized in a way similar to regularization parameter. The optimal values for gamma and C often vary depending on the dataset’s size, noise, and complexity.

GridSearchCV technique is used in our project to determine the optimal values of hyperparameters given a set of training examples. The hyperparameters 'C' and 'gamma' are given a list of numbers with kernel parameter set as 'rbf' in the parameter grid configuration. The model is also configured with a 5-fold cross-validation which uses 4 folds for training and 1 for testing (cross validation set) in an iterative manner till all the folds have been used for testing. After fitting the parameter grid on training dataset, optimal values of regularization and gamma (spread) parameter are determined and are expected to provide highest cross-validation accuracy on training data.

IV. RESULT AND OBSERVATIONS

The SVM-based IDS is evaluated using performance metrics such as accuracy, precision, recall, F1-score, and ROC-AUC, with particular emphasis on minimizing the false positive rate while maintaining high detection accuracy. After obtaining the best model by optimizing the hyperparameters, its performance is estimated against the testing set. The model has been found to achieve exceptional performance with an accuracy of 99.99% on the training set. The model shows similar results in precision, accuracy and recall as provided by the classification report.

Class 0 (Normal):

- Precision: 1.00 – Out of all predicted "Normal" instances, 100% were correctly classified.
- Recall: 1.00 – The model correctly identified 100% of the actual "Normal" instances.
- F1-Score: 1.00 – A balance between precision and recall, which indicates a perfect model.
- Support: 328816 – The total number of "Normal" instances in the test set.

Class 1 (Attack):

- Precision: 1.00 – Out of all predicted "Attack" instances, 100% were correctly classified.
- Recall: 1.00 – The model correctly identified 100% of the actual "Attack" instances.
- F1-Score: 1.00 – A perfect balance between precision and recall.
- Support: 17,214 – The total number of "Attack" instances in the test set.

Overall Metrics:

- Accuracy: 99.99% – The model correctly classified 99.99% of all instances.
- Macro Average:
 - Precision: 1.00
 - Recall: 1.00
 - F1-Score: 1.00
 - This is the unweighted mean of the metrics for both classes. Perfect precision, recall and f1 score regardless of number of samples per c
- Weighted Average:
 - Precision, Recall, and F1-Score: 1.00 – Perfect precision, recall and f1 score even when taking class imbalance into account.

RBK Kernel Accuracy: 99.99%

	precision	recall	f1-score	support
0	1.00	1.00	1.00	328816
1	1.00	1.00	1.00	17214
accuracy			1.00	346030
macro avg	1.00	1.00	1.00	346030
weighted avg	1.00	1.00	1.00	346030

Fig. 2. Classification report of the SVM model with RBF kernel showing almost perfect results in its performance with 99.99% accuracy.

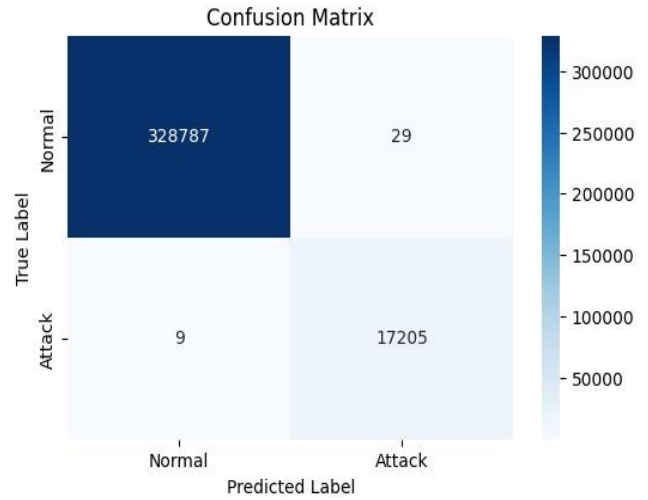


Fig. 3. Confusion Matrix of SVM model with RBF kernel showing the number of false predictions being significantly smaller than the number of true predictions.

The confusion matrix plot has shown the number of true positives and true negatives dominating the false predictions by a high margin. The model correctly identifies the majority of 'Normal' and 'Attack' instances. It has very few misclassifications indicating strong model performance.

- True Positives (Bottom-right: 17205): These are correctly predicted Attack instances. The model classified 17205 Attack instances correctly.
- True Negatives (Top-left: 328787): These are correctly predicted Normal instances. The model classified 328787 Normal instances correctly.
- False Positives (Top-right: 29): These are Normal instances that were incorrectly predicted as Attack. The model misclassified 29 Normal instances as Attack.
- False Negatives (Bottom-left: 9): These are Attack instances that were incorrectly predicted as Normal. The model misclassified 9 Attack instances as Normal.

The Receiver Operating Characteristic (ROC) curve is a plot of true positive rate or recall, ratio of true positives by the total number of actual positive instances in the y-axis against false positive rate or fallout, ratio of false positives by the total number of actual negative instances in x-axis at certain threshold values. The diagonal line in the graph is a basis for references for the current model and has an Area Under the Curve (AUC) as 0.5. Models having AUC above this random classifier perform better with the number of true instances greater than false ones. Our model shows an AUC of 1.00, indicating perfect classification such that the model can perfectly separate the classes without any error.

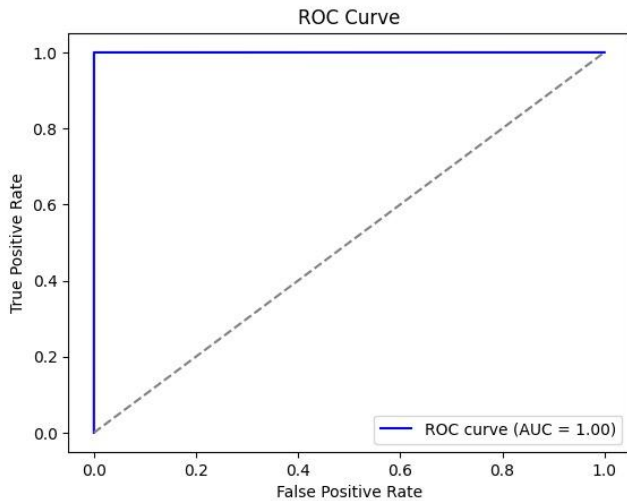


Fig. 4. ROC curve of SVM model with RBF kernel showing an AUC of 1.00 attributing to the perfect classification capability of the model.

Certain other models have been instantiated and trained with the same dataset to compare their performance with the SVM model with RBF kernel. The configurations and performance of these models are listed as follows:

- A SVM model with a polynomial kernel with degree parameter set to [2,3] which shows an accuracy of 99.99% on the same dataset.
- A RandomForestClassifier model with random configurations shows 100% accuracy.
- Simple Logistic Regression model shows 99.80% accuracy.
- XGBoost (Extra Gradient Boosting) shows 99.88% accuracy.

From the above performances, we could observe that our SVM model with RBF kernel shows a much higher accuracy than the other models, except the Random Forest Classifier. Also, the model shows great performance metrics, likely because the dataset is well-structured after the cleaning and preprocessing steps.

V. CONCLUSION

This project proposes and implements an IDS that is based on SVM algorithms to improve real-time network security by employing an anomaly-based approach which can detect previously unidentified threats due to its lack of signatures. Due to the application of a strict approach used at different stages, such as data preprocessing and model training, the system makes correct classifications of network traffic as normal and anomalous.

The evaluation using metrics like accuracy, precision, recall and F1 score proves the real-life applicability of the proposed system that has ability for accurate detection with minimal false alarms and false negative. The presented approach of an IDS as a scalable and adaptive prototype eliminates classical limitations and notably enhances threat detection, thus enhancing the general security of the network.

Throughout the development, several limitations are cited including limited computational resources, issues with data labeling and balancing, and adapting detection models to emerging network threats. These challenges are solved through various optimizations of algorithms used and the datasets selected to be implemented within the IDS, thus providing scalability and sound performance for the IDS in real-world scenarios. This structured approach has measurable results of accurately detecting emerging network threats and is versatile enough to be implemented in such dynamic network environments in contrast with any signature-based methods.

REFERENCES

- [1] P. Garcí'a-Teodoro, J. Dí'az-Verdejo, G. Macia'-Ferna'ndez and E. Va'zquez "Anomaly-based network intrusion detection: Techniques, systems and challenges", computer and security, science direct, 2009.
- [2] A. Srinivas and K. Sagar, "Anomaly Based Intrusion Detection System Using Support Vector Machine in Network Traffic," *International Journal of Creative Research Thoughts (IJCRT)*, vol. 9, no. 12, pp. 123-130, Dec. 2021.
- [3] J. Gu and S. Lu, "An effective intrusion detection approach using SVM with naïve Bayes feature embedding," *Computers & Security*, vol. 103, p. 102158, 2021, ISSN 0167-4048.
- [4] Iman Sharafaldin, Arash Habibi Lashkari, and Ali A. Ghorbani, "Toward Generating a New Intrusion Detection Dataset and Intrusion Traffic Characterization", 4th International Conference on Information Systems Security and Privacy (ICISSP), Portugal, January 2018